

CS534 - Machine Learning

Written Homework Assignment 1

Author: Vy Bui
OSUID: 934370552
Email: buivy@oregonstate.edu

Problem 1
[10pts]

(a) The log likelihood function of \mathbf{w} is

$$l(w) = \log \prod_{i=1}^N P(y_i | \mathbf{x}_i; \mathbf{w})$$

(b)

First, maximizing the log likelihood function is equivalent to minimizing the negative log likelihood function.

$$\arg \max_w \log(w) = \arg \min_w -\log(w) = \arg \min_w -\log \prod_{i=1}^N P(y_i | \mathbf{x}_i; \mathbf{w}) = \arg \min_w -\sum_{i=1}^N \log p(y_i | \mathbf{x}_i; \mathbf{w}) \quad (1)$$

Because the likelihood is Gaussian,

$$\log p(y_i | \mathbf{x}_i; \mathbf{w}) = -\frac{(y_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma^2} + \text{const} \quad (2)$$

where the const consists of all terms independent of \mathbf{w} .

Substituting (2) to (1) and dropping const results in

$$\begin{aligned} \arg \max_w \log(w) &= \arg \min_w -\sum_{i=1}^N \log p(y_i | \mathbf{x}_i; \mathbf{w}) = \arg \min_w \sum_{i=1}^N \frac{(y_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma^2} \\ &= \frac{1}{2} \sum_{i=1}^N a_i (\mathbf{w}^T \mathbf{x}_i - y_i)^2, a_i = \frac{1}{\sigma^2} \end{aligned} \quad (3)$$

(c) The gradient of $L(\mathbf{w})$ can be computed as following

$$\nabla L(w) = \frac{1}{\sigma^2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x}_i$$

Because the gradient points in the direction of steepest ascent, we need to go in the opposite direction to reach the minimum, one step at a time. Hence, the update rule is

$$\mathbf{w} \leftarrow \mathbf{w} - \gamma \nabla L(w) = \mathbf{w} - \frac{1}{\sigma^2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x}_i$$

(d) Because (3) is a quadratic function of \mathbf{w} , we can compute the global optimum by setting its gradient to 0 and solve for \mathbf{w} .

$$\begin{aligned}\frac{l}{\partial \mathbf{w}} &= \frac{1}{2\sigma^2} \frac{d}{d\mathbf{w}} ((\mathbf{y} - \mathbf{xw})^T (\mathbf{y} - \mathbf{xw})) = \frac{1}{2\sigma^2} \frac{d}{d\mathbf{w}} ((\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{xw} + \mathbf{w}^T \mathbf{x}^T \mathbf{xw})) \\ &= \frac{1}{2\sigma^2} \frac{d}{d\mathbf{w}} ((-\mathbf{y}^T \mathbf{X} + \mathbf{w}^T \mathbf{x}^T \mathbf{x}))\end{aligned}\quad (4)$$

Setting (4) to 0 results in

$$\mathbf{w}_{\text{op}}^T \mathbf{x}^T \mathbf{x} = \mathbf{y}^T \mathbf{X} \Leftrightarrow \mathbf{w}_{\text{op}}^T = \mathbf{y}^T \mathbf{x} (\mathbf{x}^T \mathbf{x})^{-1} \Leftrightarrow \mathbf{w}_{\text{op}} = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}$$

Problem 2
[10pts]

(a) Compute the log-likelihood function

$$l(w) = \log L(w) = \sum_{i=1}^N \sum_{k=1}^K \log(p(y = k | x_i)^{y_{ik}}) = \sum_{i=1}^N \sum_{k=1}^K \log\left(\frac{e^{\mathbf{w}_k^T \mathbf{x}_i}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i}}\right)^{y_{ik}}$$

(b) Compute the gradient of the log-likelihood function with regard to the weight vector \mathbf{w}_c of class c .

First, the function can be simplified as follows

$$\begin{aligned} l(w) &= \sum_{i=1}^N \sum_{k=1}^K \log\left(\frac{e^{\mathbf{w}_k^T \mathbf{x}_i}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i}}\right)^{y_{ik}} = \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log\left(\frac{e^{\mathbf{w}_k^T \mathbf{x}_i}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i}}\right) \\ &= \sum_{i=1}^N \sum_{k=1}^K y_{ik} (\mathbf{w}_k^T \mathbf{x}_i - \log \sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i}) = \sum_{i=1}^N \left(\sum_{k=1}^K y_{ik} \mathbf{w}_k^T \mathbf{x}_i - \sum_{k=1}^K (y_{ik} \log \sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i}) \right) \\ &= \sum_{i=1}^N \left(\sum_{k=1}^K y_{ik} \mathbf{w}_k^T \mathbf{x}_i - \log \sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i} \right) \end{aligned}$$

Now it is easier to calculate the gradient as follows

$$\begin{aligned} \nabla_{w_{ch}} l(w) &= \frac{\partial}{\partial w_{ch}} \left(\sum_{i=1}^N \left(\sum_{k=1}^K y_{ik} \mathbf{w}_k^T \mathbf{x}_i - \log \sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i} \right) \right) = \sum_{i=1}^N \left(y_{ic} x_i^h - \frac{x_i^h e^{\mathbf{w}_c^T \mathbf{x}_i}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i}} \right) \\ &= \sum_{i=1}^N x_i^h \left(y_{ic} - \frac{e^{\mathbf{w}_c^T \mathbf{x}_i}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i}} \right) \end{aligned}$$

where $h = 0, 1, 2, \dots, n$ and $c = 0, 1, 2, \dots, k$.

Therefore, the gradient with regard to vector w_c is

$$\nabla_{w_c} l(w) = \sum_{i=1}^N x_i (y_{ic} - P(y_i = c | x_i; w))$$

Problem 3
[10pts]

(a) Derive the posterior distribution $p(\hat{\theta}|X_1, \dots, X_n, \alpha, \beta)$ and show that it is also a Beta distribution.

First, the likelihood function can be computed as follows

$$p(X_1, \dots, X_n|\hat{\theta}) = \prod_{i=1}^n \hat{\theta}^{x_i} (1 - \hat{\theta})^{1-x_i} = \hat{\theta}^{\sum x_i} (1 - \hat{\theta})^{n - \sum x_i}$$

Then, the posterior of $\hat{\theta}|X_1, \dots, X_n$ is

$$\begin{aligned} p(\hat{\theta}|X_1, \dots, X_n, \alpha, \beta) &\propto p(X_1, \dots, X_n, \alpha, \beta|\hat{\theta})p(\hat{\theta}) \\ &= \hat{\theta}^{\sum x_i} (1 - \hat{\theta})^{n - \sum x_i} \frac{\hat{\theta}^{\alpha-1} (1 - \hat{\theta})^{\beta-1}}{B(\alpha, \beta)} = \frac{\hat{\theta}^{\alpha + \sum x_i} (1 - \hat{\theta})^{\beta + n - \sum x_i}}{B(\alpha, \beta)} \\ &= \text{Beta}(\hat{\theta}|\alpha + \sum x_i, \beta + n - \sum x_i) \end{aligned}$$

(b) For the case of observing 2 heads out of 5 tosses, the posterior distribution of θ can be calculated as follows

$$\text{Beta}(\theta|2 + 2, 2 + 5 - 2) = \text{Beta}(\theta|4, 5)$$

That for the case of observing 20 heads out of 50 tosses is

$$\text{Beta}(\theta|2 + 20, 2 + 50 - 20) = \text{Beta}(\theta|22, 32)$$

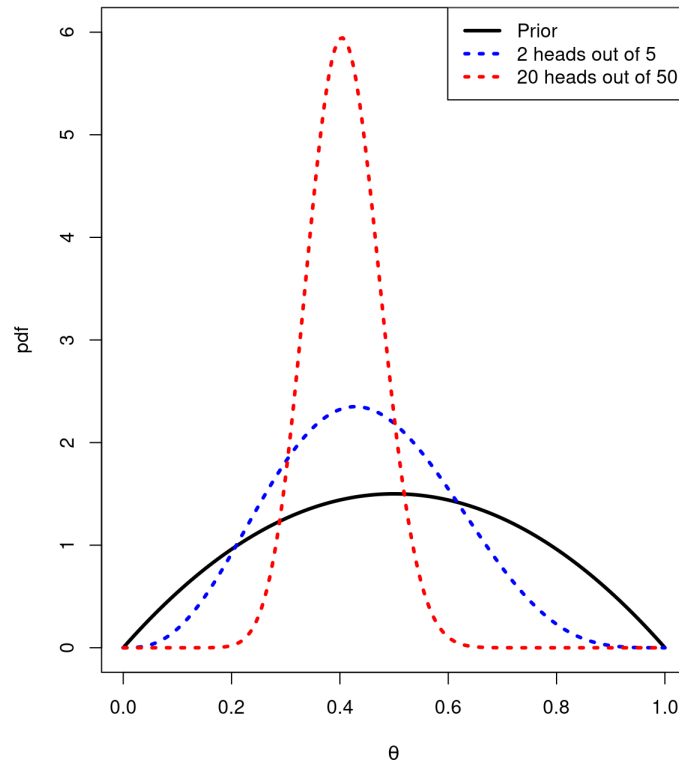


Figure 1: The posterior distributions of θ

The posterior distribution's spread will get smaller and smaller, centered at $\theta = 0.4$.